

# 7

PATENTTI- JA REKISTERIHALLITUS  
NATIONAL BOARD OF PATENTS AND REGISTRATION

Helsinki 13.12.2001

ETUOIKEUSTODISTUS  
PRIORITY DOCUMENTJC903 U.S. PRO  
10/053884  
01/22/02Hakija  
ApplicantRepublica Jyväskylä Oy  
JyväskyläPatenttihakemus nro  
Patent application no

20010136

Tekemispäivä  
Filing date

23.01.2001

Kansainvälinen luokka  
International class

G06F

Keksinnön nimitys  
Title of invention**"Menetelmä ja laitteisto tiedon uudelleenryhmittelemiseksi"**

Täten todistetaan, että oheiset asiakirjat ovat tarkkoja jäljennöksiä patentti- ja rekisterihallitukselle alkuaan annetuista selityksestä, patenttivaatimuksista, tiivistelmästä ja piirustuksista.

This is to certify that the annexed documents are true copies of the description, claims, abstract and drawings originally filed with the Finnish Patent Office.

Marketta Tehikoski  
ApulaistarkastajaMaksu 300 mk (50 € 1.1.2002 lähtien)  
Fee 300 FIM (50 EUR from 1 January 2002)

Maksu perustuu kauppa- ja teollisuusministeriön antamaan asetukseen 1027/2001 Patentti- ja rekisterihallituksen maksullisista suoritteista muutoksineen.

The fee is based on the Decree with amendments of the Ministry of Trade and Industry No. 1027/2001 concerning the chargeable services of the National Board of Patents and Registration of Finland.

Osoite:	Arkadiankatu 6 A	Puhelin:	09 6939 500	Telefax:	09 6939 5328
	P.O.Box 1160	Telephone:	+ 358 9 6939 500	Telefax:	+ 358 9 6939 5328
	FIN-00101 Helsinki, FINLAND				

## Menetelmä ja laitteisto tiedon uudelleenryhmittelyä – Metod och apparatur för att reggruppera data

5 Keksintö koskee sellaisten sääntöjen muodostamismenetelmää, jonka avulla tietoa voidaan ryhmitellä uudelleen. Lisäksi keksintö koskee järjestelyä tällaisen menetelmän toteuttamiseksi.

10 Tietoa on saatavilla valtavasti, koska sähköistä dataa kokoavia, tehokkaita laitteita on paljon ja tallennuskapasiteetti on kasvanut. Sähköiseen dataan pääsee helposti käsiksi, mutta oleellisen tiedon löytäminen onkin jo vaikeampaa. On olemassa tekniikoita, joiden avulla voidaan ennakoida ja arvioida datan sisältöä. Yleensä halutaan löytää voimassaolevaa, uutta, käyttökelpoista, ymmärrettävää tietoa. Ongelmia ilmenee esimerkiksi, jos data on epätäydellistä tai epäyhtenäistä. Data voi myös olla väliaikaista tai muuttuvaa, sitä voi olla suunnaton määrä käsiteltäväksi tai se voi olla muutakin kuin tekstimuotoista. Suuresta datamäärästä pitäisi löytää tiettyyn yritykseen tai tiettyyn ongelmaan liittyvä spesifinen tieto.

20 Eri lähteistä löytyvää tietoa on tarpeen myös yhdistää ja vertailla. Jotta epäyhtenäistä tietoa voidaan analysoida, on tieto saatettava yhtenäiseen muotoon. Lisäksi yhtenäistämistä edellyttää yritysten sisäinen ja yritysten kesken tapahtuva tiedonkulku ja -vaihto. Nykyään persoonallisuus on tärkeä erottautumiskeino markkinoilla. Persoonalliset ratkaisut kuitenkin eroavat niin sanotuista standardiratkaisuista ja vaikeuttavat siten tiedon jatkokäsittelyä. Dokumentit tulisikin esittää niin joustavasti, että ne soveltuvat ja ovat yhdistettävissä mahdollisimman moniin muihin sovelluksiin. Yleensä yksittäinen yksilöllinen ratkaisu vaatii tehokkaat laitteet hoitamaan yhteyksiä. Yhteyksistä huolehtivien laitteiden ja ohjelmistojen on oltava helppokäyttöisiä ja niiden on kasvettava yrityksen mukana. Tällaisen laitteistokokonaisuuden ylläpito ja ohjelmointi kaikki standardit ja vaatimukset huomioon ottaen on hyvin raskasta.

30 On olemassa ratkaisuja, jotka yhtenäistävät erinäisiä tiedostoformaatteja samaan, käsiteltävään muotoon. Esimerkiksi XWrap Elite (an eXtensible Wrapper Generation System Elite Version) on ohjelmistosovellus, joka pyrkii tunnistamaan toistuvat tietorakenteet HTML-dokumenteista ja luomaan niille kullekin muunnossäännöt, joilla niiden sisältö saadaan muunnettua XML-kielelle. Sovelluksen toiminta on kuitenkin sidottu HTML-rakenteeseen. XWrap Elite-ohjelmistoa on ennen tämän hakemuksen prioriteettipäivää kuvattu yksityiskohtaisemmin osoitteessa <http://www.cc.gatech.edu/projects/disl/XWRAPelite/>. Vastaavanlainen HTML-

rakenteeseen pohjautuva ratkaisu on myös W4F – World Wide Web Wrapper Factory, jota on ennen tämän hakemuksen prioriteettipäivää kuvattu tarkemmin osoitteessa <http://www.tropea-inc.com/technology/W4F/>. Lisäksi Xeroxin Euroopan tutkimuskeskuksessa on keväällä 2000 julkistettu tiedon muunnosskriptin generointimenetelmä, Wrapper Generation via Grammar Induction, joka on myös sidottu HTML-formaattiin. Tämä menetelmä on ennen tämän hakemuksen prioriteettipäivää julkaistu osoitteissa <http://turing.wins.uva.nl/~ragetli/documents/chidlovskii.ps> ja <http://turing.wins.uva.nl/~ragetli/ecml2000/ecml00a/>. Kaikki nämä mainitut, olemassa olevat ratkaisut perustuvat HTML-formaattiin, joten niiden käyttöalue on melko rajoitettu. Aiheesta on olemassa myös patentti US6151608, jossa kuvataan tiedon muuntamista tietokantataulukkoon koodia kirjoittamatta.

Keksinnön tavoitteena on tuottaa menetelmä ja järjestely, joilla ohjelmointitaidonkin käyttäjä voi muodostaa mistä hyvänsä sähköisestä tietovirrasta valitsemilleen data-alueille irrotussäännösten. Säännösten avulla etsitään ja irrotetaan data-alueita lähdedatasta.

Tavoite toteutetaan asteittain oppivan koodingenerointisovelluksen avulla.

Keksinnön mukaisesti käyttäjä saa koodia kirjoittamatta valitunlaiset data-alueet erotettua lähdedatasta jatkokäsittelyä varten.

Koodingenerointisovellus voi käyttää tiedon esittämiseen ja valintojen tekemiseen datavirrasta mitä tahansa selainta. Selaimella tarkoitetaan tässä sellaista tiedon tarkasteluvälinettä, jonka esittämään tietoon voidaan kohdistaa valintoja. Jos kyseessä on esimerkiksi HTML- tai tekstiselain, voidaan yhtenäinen jatkuva alue valita perinteisesti ”hiirellä maalaamalla”. Käyttäjä toimii selaimelle ominaisessa ympäristössä, eikä tarvitse välttämättä lainkaan koodingenerointisovelluksen tuottamaa tiedon irrottamiseen käytettävää lähdekoodia. Käyttäjällä on olemassa dataa, jotka hän jatkokäsittelyä varten haluaa eritellä tietyin osin ja/tai yhtenäistää. Nyt käyttäjän tarvitsee vain osoittaa koodingenerointiohjelmalle haluamastaan lähdedatasta kaksi tai useampia esimerkkejä. Esimerkit voidaan näyttää myös eri lähdedatoista.

Käyttäjä antaa ohjelmalle esimerkkitaupaukset, joiden perusteella ohjelma osaa irrottaa lähdedatasta tietyt rakenneosat. Lähdedatasta käyttäjä osoittaa koodingenerointiohjelmalle tietyt osat, joiden mallin mukaan ohjelma muodostaa säännösten. Säännösten perusteella koodingenerointiohjelma etsii ja irrottaa tietoa lähdedatasta. Käyttäjä voi tämän jälkeen esimerkiksi yhtenäistää epäyhtenäistä dataa samaan formaattiin. Käyttäjän antamien esimerkkien tulisi olla sisällöltään mahdollisimman

erilaisia tapauksia samantyyppisen, tunnistettavan tietorakenteen osasta. Ne voivat olla vaikkapa taulukon rivejä, luettelon alakohtia tai lomakkeita.

Koko sähköisestä tietotulvasta käyttäjä siis valitsee halutut datalähteet tai halutut osat datalähteistä ja osoittaa ne syötteenä koodingenerointiohjelmalle. Tämän tulok-  
 5 sena käyttäjä saa lähteisiin perustuvan tiedon järjestettynä sekä tuotetun irrotussään-  
 nöstön. Tulodata voidaan esittää myös halutussa formaatissa, jonka käyttäjä voi jo-  
 ko määrittää itse tai valita valmiista esimerkkimäärittäyksistä. Lisäksi ohjelma esittää  
 valitut piirteet helposti tarkasteltaviksi.

Seuraavassa selostetaan keksintöä yksityiskohtaisemmin viitaten esimerkkinä esitet-  
 10 tyyn edulliseen suoritukseen ja oheisiin kuviin, joissa

kuva 1 esittää keksinnön peruseriaa-

kuva 2 esittää erään edullisen sovellusmuodon mukaista muunnoskriptin luomis-  
 ta ja

kuva 3 esittää erään edullisen sovellusmuodon mukaista tokenisoitujen esimerkki-  
 15 en käsittelyä.

Hajanaisia, eri muodoissa olevia dokumentteja ja tiedostoja on usein tarpeen yhten-  
 näistää, jotta niitä voidaan jatkokäsitellä. Yhtenäistäminen on tärkeää yritysten väli-  
 siä asioita käsiteltäessä, koska on käytännössä mahdotonta vaatia, että kaikki yhteis-  
 työ- tai asiakassuhteissa toimivat yritykset yhtenäistävillä tietojärjestelmänsä ja  
 20 ohjelmistonsa. Yhtenäistettävyyden käsittelyssä helpottaa paitsi yritysten sis-  
 äistä tiedonkäsittelyä, erityisesti julkaisu- ja yhteistyöverkostojen toimintaa. On  
 oleellista, että sähköisestä informaatiotulvasta voidaan valita itselle tärkeät alueet ja  
 että näitä osioita voidaan tarkastella, vertailla ja käsitellä edelleen yhtenäisenä ko-  
 konaisuutena.

Kuvassa 1 on esitetty keksintö pääpiirteissään. Koodingenerointisovellus on datan  
 erotteluun, järjestämiseen ja yhtenäistämiseen sovitettu työkalu. Koodingenerointi-  
 sovellus käsittää kolme osaa: koodingenerointikomponentin 102, muodostettavat ir-  
 rotussäännöt 103 sekä irrotuskomponentin 104. Lähdemateriaali 101 voi olla mitä  
 25 hyvänsä tietoteknistä materiaalia, kuten esimerkiksi tiedostoja, dokumentteja tai jat-  
 kuvaa tietovirtaa. Ainoa edellytys lähdeformaatille on, että välittömästi jokaisen  
 30 automaattisesti irrotettavaksi halutun tietokentän edellä tai jälkeen on oltava vähin-  
 tään yhden käsiteltävän datan perusyksikön eli tokenin mittainen kenttäerotin, joka  
 toistuu samanlaisena kaikissa tietueissa juuri ennen kyseistä tietokenttää tai välitti-  
 mästi sen jälkeen. Itse irrotettava tietue voi kuitenkin sisältää erilaisia kenttäerotti-

mia. Käyttäjä valitsee koko lähdemateriaalin tai tiettyjä lähdemateriaalin osia, data-alueita. Näiden valittujen data-alueiden tietorakenteen tulee olla litteä eli tietoraken-  
teessa ei saa esiintyä eri hierarkiatasoja. Tällainen litteä tietorakenne on esimerkiksi  
5 taulukolla tai listalla. Jos käyttäjä haluaa käsitellä vain tiettyä osaa tietorakenneyksi-  
köstä, vaikkapa taulukkoa, hänen tulee antaa koodingenerointisovelluksen koodin-  
generointikomponentille 102 esimerkkinä vähintään kaksi sellaista taulukon riviä,  
jotka eroavat toisistaan mahdollisimman paljon. Mitä enemmän esimerkit poikke-  
vat toisistaan, sitä vähemmän esimerkkejä tarvitaan toivotunlaisen irrotustuloksen  
aikaansaamiseksi. Valittujen esimerkkirivien mahdollisimman monessa sarakkeessa  
10 tulee siis olla eri tietoa.

Koodingenerointisovelluksen irrotussäännöt tuottava komponentti 103 muodostaa  
esimerkkien perusteella säännöstön, jonka perusteella irrotuskomponentti 104 etsii  
ja irrottaa halutunlaisia data-alueita lähdemateriaalista 101. Varsinainen irrotuskom-  
ponentti voi olla mikä hyvänsä asiaan soveltuva komponentti. Irrotuksen tuloksena  
15 saadut data-alueet voidaan jatkokäsitellä halutulla tavalla, esimerkiksi tallentaa. Mi-  
käli ne tallennetaan tiettyyn yhteen kohteeseen, käyttäjä saa dokumentin 105, joka  
voi olla käytännössä millainen hyvänsä riippuen käyttäjän valinnoista ja alkuperäi-  
sistä lähdedatoista. Alkuperäinen dokumentti säilyy muuttumattomana. Tuotettu tu-  
los sisältää käyttäjän syöttämien esimerkkien mukaiset data-alueet, esimerkiksi  
20 kaikki taulukon rivit, käyttäjän määrittämässä muodossa ja järjestyksessä. Käyttäjä  
voi määrittää itse esimerkiksi kohdeformaatin muutaman parametrin avulla tai valita  
valmiista, olemassa olevista formateista. Jos käyttäjä valitsee valmiin formaatin,  
formaatti määritetään esimääriteltujen parametrien avulla. Saatavaa kohdedataa  
varten käyttäjä voi määritellä myös kenttien järjestyksen ja/tai jättää haluamansa  
25 tietokentät kokonaan pois käytöstä.

Sitä osaa lähdedatasta, joka valitaan muunnettavaksi, kutsutaan data-alueeksi.  
Tyypillisesti data-alue on jokin alue, jonka sisällä toistuu vain yksi tietorakenne.  
Data-alue voi olla esimerkiksi taulukko, joka siis sisältää taulukossa toistuvat rivit.  
Data-alue voi myös olla esimerkiksi kappale tekstiä, otsikko tai luettelo. Data-alueet  
30 voivat olla keskenään täysin erilaisia, mutta kunkin data-alueen sisällä toistuu yksi  
säännöllinen tietorakenne. Tämä toistuva tietorakenne on helposti osoitettavissa  
koodingenerointisovellukselle ja edelleen koodingenerointisovelluksen löydettävissä  
ja erotettavissa lähdedatasta. Taulukossa tällainen toistuva ominaisuus olisi esimer-  
kiksi tietosisällön järjestäytyminen riveihin ja/tai sarakkeisiin. Tekstikappaleita  
35 puolestaan ryhmittelee yleensä niiden alussa ja lopussa olevat tyhjät tilat ennen seu-  
raavan osion alkua. Otsikko on varustettu tietyillä tunnistettavilla ominaisuuksilla,

kuten esimerkiksi lihavoinnilla ja/tai tekstin sijoittelulla. Luettelo sisältää yleisesti jonkin tunnistesymbolin, esimerkiksi viivan tai laatikon ennen varsinaista asiasisältöä ja rivinvaihdon sen jälkeen. Oleellista siis on, että on olemassa jokin tietyille osioille yhteinen tietorakenne, jonka koodingenerointisovellus löytää käyttäjän antamien esimerkkien avulla.

- Data-alueita voidaan jakaa osiin myös sisällön perusteella. Näitä vaihtelevan tai kiinteän pituisia yksiköitä kutsutaan tokeneiksi. Menetelmää, jonka avulla merkkijonoja jaetaan osiin sisällön perusteella, kutsutaan tokenisoinniksi. Yksi tokeni sisältää tyypimerkinnän eli nimen sekä itse datasisällön. Oletuksena on, että yksi tokeni eli käsittelyn yksikkö on yksittäinen merkkijärjestelmän symboli, kuten kirjain, numero tai välimerkki. Usein osoittautuu kuitenkin kannattavammaksi määritellä tokenit siten, että yksi tokeni sisältää tietyn datassa toistuvan merkkijonon. Tokenityypit ovat siis eräänlaisia yleistyksiä. Esimerkiksi "numero"-tyyppisiä tokeneja ovat kaikki numerot, jos tokenityypin säännöllinen lauseke on "[0-9]" tai jos "sana"-tyyppisen tokenin säännöllinen lauseke on "[A-Za-z]+"
- 15 tokeni voi sisältää yhden tai useamman peräkkäisen merkin väliltä A – Z tai a – z. On myös mahdollista yksilöidä, että esimerkiksi vain ja ainoastaan merkkijono "B2B" on "test2"-tyyppinen tokeni. Tokenisoidut merkkijonot ovat helpommin hallittavissa kuin yksittäiset symbolit.
- 20 Tokenisointi helpottaa ja nopeuttaa datan käsittelyä. Käytetyn tokenisoinnin ansiosta nyt kyseessä oleva sovellus ei ole sidottu mihinkään tiettyyn formaattiin, vaan riittää, että irrotettavaksi valitun tietokentän edellä tai sen jälkeen on vähintään yhden tokenin mittainen kenttäerotin. Tokenisointi siis monipuolistaa ja nopeuttaa koodingenerointisovelluksen käyttöä, sillä vähäinenkin pohjatieto käsiteltävästä
- 25 formaatista tokenisointimäärittäykseksi kirjoitettuna mahdollistaa parempiin irrotussääntöihin pääsemisen pienemmillä esimerkkimäärillä. Vastaavasti virheellisen tokenisoinnin käyttö voi johtaa siihen, ettei algoritmi pysty tuottamaan ollenkaan toimivia irrotussääntöjä. Tokenisoimalla data ja suorittamalla tämän jälkeen kaikki käsittelyoperaatiot tokenisoituun dataan voidaan varmistaa koodingenerointisovelluksen oikeanlainen toiminta. Tokenisointia hyväksikäyttäen saavutetaan myös kulloiseenkin tilanteeseen sopiva tarkkuus lähdemateriaalin tarkastelulle. Tässä keksinnössä tokenisointi on apuväline irrotussääntöjen generoinnissa, mutta sääntöjä tulkitseva tiedon irrotuskomponentti ei enää tiedä tokenisoinnista mitään, eikä käytä sitä mihinkään.
- 30
- 35 Käyttäjä osoittaa koodingenerointisovellukselle esimerkit. Esimerkkien osoittamisella tarkoitetaan tässä, että käyttäjä osoittaa esimerkit olemassa olevasta datasta ja/tai

tuottaa esimerkit itse ja/tai muokkaa olemassa olevia esimerkkitapauksia korostaakseen niiden yhteneviä osia.

Käyttäjän koodingenerointisovellukselle antamien esimerkkien tulisi olla mahdollisimman erisisältöiset ei-yhteneviltä osiltaan, jotta koodingenerointisovellus yksiselitteisesti löytäisi esimerkkitapausten yhtenäiset osat, joita se käyttää luodessaan datanirrotussäännöstöä. Koska esimerkkien on oltava mahdollisimman erilaiset sisällöltään, tokenisoinnin avulla mahdollistetaan tietynlaisten datalähteiden käyttö. Esimerkiksi HTML-tiedosto voidaan jakaa tageihin ja niiden välisiin tekstipaloihin. Jos HTML-formaattia käsiteltäisiin ilman tokenisointia, eli oletusarvoisesti yksi symboli muodostaisi yhden tokenin, esimerkkien osoittaminen koodingenerointisovellukselle olisi työlästä, jollei mahdotonta. Käyttäjältä kuluisi kohtuuttomasti aikaa sopivien esimerkkien löytämiseen, koska näiden sisältöjen tulisi olla symbolitasolla täysin erilaiset. Esimerkkitiedon yksilöimiseksi tarvittaisiin siis hurjan paljon esimerkkejä tai sitten niiden laatua tulisi voida oleellisesti parantaa. Tämä laadun parannus tehdään nimenomaan tokenisoimalla: kun HTML-formaatin mukainen data on tokenisoitu tageiksi ja niiden väliin jääviksi merkkijonoiksi, yhdenkin yksittäisen symbolin ero merkkijonojen sisällä riittää ilmaisemaan koodingenerointisovellukselle, että alue halutaan irrottaa ja säilyttää yhtenäisenä.

Tässä tekstissä yksittäisestä merkkijärjestelmän merkistä käytetään nimitystä ”symboli”. Useammasta peräkkäisestä symbolista käytetään kuitenkin nimitystä ”merkkijono”. Mainitut irrotettavien data-alueiden alku- ja loppumerkit eivät ole yksittäisiä symboleja, vaan merkkijonoja tai säännöllisiä lausekkeita. Jokerimerkkikään ei välttämättä ole vain yksi merkki (\*), vaan sitä voi vastata jokin merkkijono (esimerkiksi .\*) järjestelmästä, ohjelmointikielestä ja sovelluksesta riippuen.

Koodingenerointisovelluksen käyttäjä voi halutessaan itse määritellä tokenisointisäännöt. Yleensä tokenisointisäännöt valitaan valmiista asetuksista, joista esimerkkinä HTML-kieleen sopivasta tagien ja niiden väliin jäävien merkkijonojen erittely. Käyttäjä voi myös muokata valmiita tokenisointisääntöjä.

Esimerkkinä esimerkkitapausten erilaisuudesta voidaan tarkastella tilannetta, jossa esimerkkitapauksia on neljä ja kukin niistä koostuu kolmesta peräkkäisestä kentästä. Se, että kussakin esimerkkitapauksessa on kolme kenttää, tarkoittaa, että esimerkkitapauksilla on samanlainen rakenne. Sisällön erilaisuus tarkoittaa, että jos tarkastellaan yhdessä kaikkien esimerkkitapausten ensimmäisiä (tai toisia, tai kolmansia) kenttiä, ne eivät ole kaikissa esimerkkitapauksissa samat. On kuitenkin sallittua, että vaikkapa ensimmäisen esimerkkitapauksen ensimmäinen kenttä on sama kuin toisen

ja kolmannen esimerkkitapauksen ensimmäinen kenttä, kunhan neljännen esimerkkitapauksen ensimmäinen kenttä on sisällöltään erilainen.

Koodingenerointisovellus tallentaa käyttäjän syöttämät esimerkit, joten esimerkkejä voidaan lisätä yksi kerrallaan valiten niitä eri lähteistä. Tyypillisesti myös irrotussäännöt tallennetaan, jolloin samoja, tallennettuja sääntöjä voidaan käyttää myöhemminkin tarvitsematta generoida niitä uudelleen.

Käyttäjän lähdedatasta valitsema esimerkki on merkkijono. Tällaisen irrotettavan merkkijonon edellä tai sen jälkeen on oltava vähintään yhden tokenin mittainen kenttäerotin, joka toistuu samanlaisena kaikissa tietueissa. Tietueen sisällä kenttäerottimien ei kuitenkaan tarvitse olla samoja. Tietorakenteen perusteella irrotettavissa olevat kentät paikallistetaan ja yksilöidään, eli kentille annetaan järjestysnumero. Tietokenttien rajat määräytyvät kulloinkin käytettävän tokenisoinnin mukaisesti. Käyttäjän antamien esimerkkien perusteella koodingenerointisovellus muodostaa editointiskriptin (edit script). Editointiskriptissä luetellaan ne muutokset, joilla ensimmäisen syötetyn esimerkkitapauksen sisältö saadaan samaksi kuin toisen esimerkkitapauksen sisältö. Editointiskriptissä toistuva perusyksikkö on muutostieto. Tämä muutostieto sisältää indeksiosoitteet molempiin merkkijonoihin sekä korvattavan ja korvaavan tiedon kokonaisuudessaan. Kohdistamalla määrätyllä tavalla editointiskriptien muutokset pisimpään syötettyyn esimerkkilausekkeeseen, jota kutsutaan säännölliseksi lausekkeeksi, muodostetaan tuloksena saadun säännöllisen lausekkeen avulla irrotusskripti, joka sitten etsii lähdedatasta merkkijonoa syötetietojen perusteella. Edelleen irrotusskriptin mukaiset tietyt, valitut data-alueet irrotetaan lähdedatasta jatkokäsittelyä varten.

Vaikka esimerkit koodingenerointisovellukselle näytettäisiin tiedostoista, varsinainen tiedon irrotus generoitujen sääntöjen avulla voidaan tehdä tiedostojen lisäksi myös esimerkiksi jatkuvasta, päättymättömästä tai tuntemattoman pituisesta tietovirrasta. Irrotettavan data-alueen pituus tiedetään siis vasta siinä vaiheessa, kun tietovirrasta tai erillisessä kontrollivirrasta löydetään jokin ennaltamäärätty lopetusmerkki.

Kun käyttäjä on valinnut esimerkkitapaukset, hän voi testata yksittäisen data-alueen irrotussääntöjen toimintaa esikatselu-tilassa koodingenerointisovelluksen selaimessa. Jos käyttäjä havaitsee, että ennen hänen esimerkiksi tarkoittamaansa data-aluetta on jotain ylimääräistä dataa, käyttäjä voi määritellä data-alueen alku- ja loppumerkit tai yrittää täsmentää olemassa olevia. Käyttäjän tulee valita ennen toivottua irrotettavaa data-aluetta edeltävä merkkijono siten, että se olisi mahdollisimman yksiselit-



teinen ja täsmällinen kertomaan sovellukselle, että juuri tällaisen tietyn merkkijonon jälkeen tulevaa dataa lähdemateriaalista halutaan irrottaa. Näin käyttäjä voi esikatse-  
lu-tilassa tarkistaa, että valittujen esimerkkitapausten tyyppistä tietoa irrotetaan oi-  
keilta alueilta. Vastaavasti voidaan varmistaa sovelluksen toiminta data-alueen lop-  
pumerkin osalta.

Esimerkit voidaan myös syöttää täysin erillisinä vaikkapa suoraan näppäimistöltä. Kertaalleen osoitettuja esimerkkejä on mahdollista muuttaa: jos käyttäjä huomaa yhden kokeilun jälkeen, ettei kahden esimerkin välillä ollut riittävää sisällöllistä eroa, hän voi lisätä eroja syöttämällä ainakin toiseen esimerkkiin joitain sellaisia symboleja, joita toisessa ei ole.

Koodingenerointisovelluksen selaimen ominaisuuksista riippuu, näkyvätkö data-alueen alku- ja loppumerkkejä kuvaavat merkkijonot käyttäjälle suoraan lähdedatas-  
ta ja/vai korostetaanko niiden välinen alue selkeästi havaittavaksi. Tarkoitushan on, että lähdemateriaalin alusta lähtien ensimmäinen löydetty alkua tarkoittava merkki-  
jono on irrotettavaksi haluttua data-aluetta edeltävä kohta. Alkumerkistä eteenpäin lähdemateriaalista etsitään loppumerkkiä, johon irrotettava data-alue päättyy. Jos data-alueen alku- tai loppumerkki puuttuu eli on tyhjä, koodingenerointisovellus tulkitsee irrotettavan data-alueen aluksi lähdemateriaalin alun ja vastaavasti irrotet-  
tavan data-alueen lopuksi lähdemateriaalin lopun. Data-alue-määrittäminen on siis aina olemassa, eikä se ole millään tavoin riippuvainen lähdemateriaalin formaatista, ra-  
kenteesta tai sisällöstä. Käyttäjä voi myös muokata tiedon irrotussäännösten ohjel-  
makielistä koodia, jos esimerkiksi sopivia esimerkkejä ei löydy riittävästi.

Jos käyttäjä haluaa tuloksena tiedon yhtenäisessä muodossa, käyttäjä voi määrittellä kohdeformaatin muutamalla parametrilla tai valita jonkin valmiin, kuten esimerkiksi XHTML-  
taulukon, pohjaksi. Halutessaan käyttäjä voi siis määrittää, missä muodossa ja millä lisämerkinnöillä varustettuina irrotetut tietueet kirjoitetaan kohteeseen. Jos tuloksena saatavaa kohdedataa halutaan käsitellä esimerkiksi XML-selaimella, voidaan valmiiksi määrittää tiettyä lähdemateriaalista irrotettua osaa ympäröivät tagit, jotta tuleva kohdedata on heti havainnollisemmassa ja helpommin käsiteltäväs-  
sä muodossa. Seuraavassa on esimerkkinä XML-formaatin mukaisia tageja ja niiden välissä se lähdemateriaalin osa, jonka tagit ympäröivät:

```

<datastream>koko kohdedata</datastream>
<dataarea>data-alue</dataarea>
<record>tietue</record>
<field>kenttä</field>
5  <address>kenttä</address>
    <price>kenttä</price>.

```

Data-alueella tarkoitetaan siis esimerkiksi tauluja tai taulukoita, joista kukin sisältää yhden tai useampia tietueita eli taulun/taulukon rivejä, joista kukin sisältää yhden tai useampia kenttiä eli taulun/taulukon rivien soluja/sarakkeita.

- 10 Keksintö ei sinänsä edellytä, että irrotetut tietueet tallennetaan johonkin tiettyyn yhteen kohteeseen yhtenäisessä formaatissa. Kun tieto on irrotettu lähteestä, sitä voidaan jatkokäsittää monin tavoin. Esimerkkejä mahdollisesta jatkokäsittelystä ovat irrotetun tiedon antaminen jollekin toiselle ohjelmistokomponentille tai laitteelle, irrotetun tiedon lähettäminen eteenpäin esimerkiksi verkkoyhteyttä pitkin ja irrotetun tiedon osittainen tai täydellinen hävittäminen.
- 15 Irrotetut data-alueet ja/tai niiden osat (tietueet, kentät) voidaan irrotuksen jälkeen järjestää uudelleen erilaisiin järjestyksiin joko ennen niiden jatkokäsittelyä tai osana jatkokäsittelyä.

- Tarkastellaan yksityiskohtaisemmin koodingenerointisovelluksen toimintaa ja muunnoskriptin tuottamista yhden data-alueen osalta kuvan 2 avulla. Koodingenerointisovellukselle osoitetaan syötteenä lista data-aluemäärittelyksiä, joista jokainen sisältää vähintään kaksi esimerkkiä ja merkkijonon tai säännöllisen lausekkeen valitun data-alueen alkamis- ja päättymiskohdan tunnistamiseksi. Tietorakenne-esimerkkien sisälle ei tarvitse erikseen laittaa mitään aloitus- tai lopetusmerkkejä. Aluksi kohdassa 201 tokenisoidaan kaikki käyttäjän syöttämällä esimerkkilistalla olevat esimerkit.
- 20 Kuvan algoritmi ei tarvitse data-alueiden alku- ja loppumäärittelyksiä esimerkkien käsittelyyn, vaan se yksinkertaisesti lisää data-alueen rajaamiseksi tarvittavat, valmiina saamansa tiedot irrotusskriptiin.

- Seuraavaksi 202 esimerkit järjestetään tokenimääräisesti laskevaan pituusjärjestykseen. Kohdassa 203 listan ensimmäinen ja siis pisin esimerkki kopioidaan ja se määritetään säännölliseksi lausekkeeksi R. Tämän jälkeen käsitellään listan seuraavaa esimerkkilauseetta 204. Seuraavaa, toiseksi pisintä lausetta merkitään esimerkkilausekkeella E. Kohdassa 205 verrataan säännöllistä lauseketta R esimerkkilausekkeeseen E käyttämällä sinänsä tunnettua Diff-vertailualgoritmia, jonka tarkempi kuvaus ennen tämän hakemuksen prioriteettipäivää löytyy osoitteesta
- 30 <http://www.cs.arizona.edu/people/gene/PAPERS/diff.ps>. Tämä vertailualgoritmi

palauttaa lyhimmän mahdollisen editointiskriptin D eli säännöllisen lausekkeen R ja esimerkkilausekkeen E välisen muutostiedon. Jos toteutetaan editointiskriptin D sisältämät muutokset säännölliselle lausekkeelle R, saadaan esimerkkilause E.

5 Keksintö ei sinänsä edellytä nimenomaan tunnetun Diff-algoritmin käyttöä, vaan sen tilalla voidaan käyttää mitä tahansa sellaista vertailevaa algoritmia, joka palauttaa editointiskriptiä vastaavat ohjeet siitä, miten tietystä kohteesta saadaan lyhimmillä mahdollisella tavalla muuntamalla tietty toinen kohde. Jos on olemassa vähintään kaksi lyhintä mahdollista tapaa, keksinnön kannalta ei ole väliä sillä, mikä näistä tavoista valitaan. Lyhin mahdollinen tapa tarkoittaa, että ns. editointietäisyys eli tarvittavien muutosten lukumäärä muunnettaessa lähdettä kohteeksi on mahdollisimman pieni.

15 Seuraavaksi 206 luetaan Diff-vertailualgoritmin palauttamasta editointiskriptistä D muutosoperaatiot yksi kerrallaan. Editointiskriptiä D ja säännöllisen lausekkeen R sisältämiä tokeneja vertaillaan keskenään 207. Säännöllistä lauseketta R muutetaan siten, että ne tokenit, joihin D:stä luettu muutostieto viittaa, merkitään jokerimerkillä, jolloin jäljelle jäävät alkuperäisinä vain R:n muuttumattomat tokenit. Seuraavaksi tutkitaan, onko kaikki D:n muutosoperaatiot jo läpikäyty. Jos muutosoperaatioita on vielä lukematta, palataan kohtaan 206. Jälleen tokenit, joihin muutostieto kohdistuu, merkitään säännöllisessä lausekkeessa R jollain tietyllä jokerimerkillä (\*-wildcard).

20 Kohdassa 209 tarkastetaan, onko esimerkkilistalla vielä lisää esimerkkilausekkeita. Kaikki listan esimerkit käydään yksi kerrallaan läpi eli jos lisäesimerkkejä löytyy, jatketaan kohdasta 204. Mikäli kaikki esimerkkilausekkeet on jo käsitelty, edetään kohtaan 210 ja poistetaan tuotetusta säännöllisestä lausekkeesta peräkkäiset jokerimerkit siten, että vain ensimmäinen niistä jää jäljelle. Seuraavaksi säännöllinen lauseke R muunnetaan osaksi muunnosskriptiä 211. Alun perin (ennen tätä vaihetta) muunnosskripti on tyhjä. Tuotettava muunnosskripti etsii lähdemateriaalista merkkijonoja syöttötietojen perusteella. Kun säännöllisessä lausekkeessa on jokerimerkki, lähdemateriaalista etsitään se merkkijono, joka säännöllisessä lausekkeessa on jokerimerkin jälkeen. Jokerimerkki tavallaan korvaa pienimmän mahdollisen määrän

25 mitä hyvänsä lähdemateriaalin merkkejä siihen merkkijonoon asti, joka säännöllisessä lausekkeessa esiintyy jokerimerkin jälkeen. Tällaista jokerimerkin tulkintaa nimitetään usein reluctant-tulkinnaksi.

35 Muunnosskriptin eteen lisätään vielä kohdan 212 mukaisesti koodi, jolla yksi data-alue kokonaisuutena irrotetaan lähdedatasta. Tämän lähdedatan sisältöä koodingene-rointisovellus tässä esimerkissä käsittelee ja muuntaa halutut osat siitä kohdefor-

maatin mukaisiksi. Tämän jälkeen muunnoskriptiin voidaan vielä lisätä tämän esimerkin mukaisesti käyttäjän mahdollisesti antamat tai ennaltamääritetyt parametrit, jotka määrittävät kohdeformaatin. Koodingenerointisovellus siis tuottaa muunnoskriptin, joka irrottaa kaiken data-alueilla esiintyvän annettuihin tietorakenne-esimerkkeihin täsmäävän tiedon jatkokäsittelyä varten, joka tässä esimerkkitapauksessa on konvertointi yhtenäiseen kohdeformaattiin.

Eräs vaihtoehto muunnoskriptin automaattiselle muodostamiselle kohtien 211 ja 212 mukaisesti on jatkaa kohdasta 210 siten, että muodostettu säännöllinen lauseke R otetaan talteen ja käsitellään muulla tavoin kuin liittämällä se osaksi muunnoskriptiä. Se voidaan esimerkiksi käyttää ”manuaalisesti” uudelleen jonkin toisen ohjelmatuotteen osana tai sitä voidaan soveltaa pelkkänä irrotusskriptinä tekemättä konvertointia kohdeformaattiin.

Tarkastellaan vielä kuvan 3 avulla säännöllisen lausekkeen muodostamista. Tässä esimerkkitapauksessa käyttäjä on osoittanut koodingenerointisovellukselle kolme esimerkkilauseketta, jotka on tokenisoitu. Kaksi pisintä esimerkkilauseketta koostuu 8 tokenista ja lyhin esimerkkilauseke 6 tokenista. Säännölliseksi lausekkeeksi R 301 merkitään pisin esimerkkilauseke eli tässä tapauksessa kumpi hyvänsä pidemmistä 8 tokenin esimerkkilausekkeista. Jäljelle jäänyt pisimmistä esimerkkilausekkeista merkitään E:ksi 302. Lausekkeet R ja E osoitetaan parametreina tunnetulle Diff-vertailualgoritmillemme, joka palauttaa editointiskriptin D 303. Nyt editointiskriptin muutostiedot on esitetty siten, että kunkin rivin ensimmäiset numerot viittaavat säännöllisen lausekkeen R 301 korvattaviin tokeneihin ja seuraavat numerot esimerkkilausekkeen E 302 korvaaviin tokeneihin. Lisäksi editointiskripti D 303 sisältää niiden tokenien sisällön, joissa muutoksia tapahtuu. Kuvan 3 editointiskriptin D 303 mukaisesti säännöllisen lausekkeen R 301 tokeni R[2] muutetaan vastaamaan esimerkkilausekkeen E 302 tokeneja E[2] ja E[3]. Säännöllisen lausekkeen R 301 tokenit R[5] ja R[6] on korvattava esimerkkilausekkeen E 302 tokenilla E[6]. Kaikki muutuvat tokenit on kuvassa ilmaistu katkoviivanuolella ja samana pysyvät, toisiaan vastaavat tokenit on kuvattu yhtenäisellä nuolella.

Seuraavaksi muodostetaan uusi säännöllinen lauseke R' 304 siten, että muuttumattomat tokenit säilyvät säännöllisessä lausekkeessa R 301 ennallaan ja kaikki esimerkkilausekkeesta E 302 eroavat, editointiskriptin mukaisen muutostiedon mukaan muuttuvat, katkoviivanuolella merkityt tokenit korvataan jokerimerkillä, joka tässä esimerkissä on \*. Uudeksi esimerkkilausekkeeksi E' 305 merkitään seuraava käyttäjän syöttämä esimerkki, minkä jälkeen tuotetaan R' 304 ja E' 305 välinen muutostieto D' 306 Diff-vertailualgoritmin avulla. Diff-vertailualgoritmin palauttamasta edi-

tointiskriptistä D' 306 luetut muutokset on merkitty katkoviivanuolilla ja R':n 304 ja E':n 305 muuttumattomat tokenit on yhdistetty yhtenäisellä nuolella. Lisäksi editointiskriptin toinen muutostieto koskee tokenia R[4], jolla ei ole vastinetta esimerkkilausekkeessa E' 305. Tällöin tokeni R[4] korvataan tyhjällä ja seuraavaa säännöllistä lauseketta muodostettaessa sitä käsitellään kuten muitakin korvattavia tokeneja. Muodostetussa säännöllisessä lausekkeessa R'' 307 ovat jäljellä muuttumattomat tokenit R[1], R[3], R[7] ja R[8]. Kaikki muuttuneet tokenit on korvattu valitulla jokerimerkillä \*. Tämän jälkeen säännöllisestä lausekkeesta R'' 307 poistetaan peräkkäiset jokerimerkit (ei esitetty kuvassa) siten, että vain ensimmäiset niistä, jäävät jäljelle. Esimerkki lausekkeessa R'' 307 on nyt neljä kenttäerotintokenia, joiden väliin muodostuu kaksi eri kenttää. Ensimmäinen varsinaista tietosisältöä sisältävä kenttä on tokenien R[1] ja R[3] välissä ja toinen kenttä sijaitsee tokenien R[3] ja R[7] välissä. Kenttäerottimet R[7] ja R[8] esiintyvät peräkkäin ja muodostavat tämän säännöllisen lausekkeen loppumerkin. Tämän jälkeen säännöllinen lauseke R'' 307 on valmis ja voidaan esimerkiksi lisätä osaksi irrotusskriptiä.

Mikäli jossakin esimerkkilausekkeessa ei ole mitään rakennetta, edellä selostettu toiminta johtaa siihen, että säännöllinen lauseke R'' on lopulta pelkkä yksi jokerimerkki. Tällöin tuloksena syntyvä muunnoskripti irrottaa koko sen data-alueen, johon irrotussäännöstöä yritetään soveltaa. Irrotuksen tuloksena kaikki irrotettu tieto sijaitsee yhden irrotetun data-alueen yhden tietueen yhdessä kentässä.

Koodingenerontisovelluksen koodingenerointikomponentti (102) tuottaa koodikielisen irrotusskriptin, jonka pohjalta varsinainen irrotuskomponentti suorittaa valikoidun tietojen irrotusta sisään tulevasta datasta. Edullisesti irrotusskripti tuotetaan jollain tunnetulla koodikielellä sovelluksesta riippuen. Edullisesti koodingenerontisovellus toimii kaikissa Java 1.2.2 tai sitä uudemmissa versioissa. Tulosteena saadaan irrotetut tiedot koodikielisen skriptin määrittämässä muodossa ja järjestyksessä. Tässä esimerkinomaisesti mainittu DEL (Data Extraction Language) -kieli on tehty nimenomaan ja ainoastaan tiedon irrotusta varten, mutta se ei millään tavoin rajoita keksintöä käyttämään vain tätä tiettyä kieltä, vaan edullisesti voidaan käyttää esimerkiksi Perl-, Python-, REBOL- tai OmniMark-suorittimia. Kieli toimii formaatin konvertoijan tavoin eli se ei irrota esimerkiksi lauseen subjektia, vaan toimii symbolitasolla. Säännöllisiin lausekkeisiin pohjautuvien irrotussääntöjen muokkaamiseen on oma editorinsa, joka riippuu kulloinkin tiedon irrotukseen käytettävästä komponentista. Irrotettava tietokenttä on merkkijono, joka voidaan skriptiä muokkaamalla asettaa siirrettäväksi haluttuun kohtaan koodingenerontisovelluksen irrotuskomponentin (104) tuottamaa tulosdataa.

Koodingenerointisovellus voi käyttää tiedon näyttämiseen ja tiedosta valintojen tekemiseen lähes mitä hyvänsä selainta. Selaimelta edellytetään, että sen datan lataamiseen (tiedoston/tietovirran avaaminen/sulkeminen), esittämiseen ja valitsemiseen liittyvä toiminnallisuus on jonkinlaisen ohjelmointirajapinnan kautta ohjattavissa ja  
 5 säädeltävissä, ja että kaikki tiedot käyttäjän tekemistä valinnoista saadaan heti tapahtumahetkellä koodingenerointisovelluksen käytettäväksi. Käytetystä selaimesta riippuu, miten ja millaisia osia jatkuvasta tietovirrasta esitetään käyttäjälle. On myös mahdollista käyttää useampia selaimia yhtäaikaaisesti siten, että esimerkiksi kulloinkin käsiteltävä lähdedata sijaitsee omassa selaimessaan, valitut esimerkit ovat tallettuna toisessa selaimessa ja tiedon irrotuksen tulosten esikatselu käyttää kolmatta selainta.  
 10

Tiedon irrotussääntöjen generointi sekä tiedon irrotus ja konversio haluttuun kohdeformaattiin haluttaessa voidaan toteuttaa palveluna verkon yli. Tällöin joko koodingenerointikomponentti, irrotuskomponentti tai molemmat voivat sijaita palvelimella  
 15 muun sovelluksen toimiessa käyttäjän työasemalla.

Sovelluksen suorituksen tuloksena saadaan siis halutun sisältöinen tietopaketti. Käyttäjän ei tarvitse osata ohjelmoida tai muuntaa sovellusta, vaan hän saa esimerkkiensä mukaisen tuloksen mistä hyvänsä haluamastaan lähdedatasta. Saatuja tietoja voidaan edelleen jalostaa, tilastoida ja julkaista yhtäaikaisesti eri muodossa  
 20 käyttäen ulkoisia sovellusohjelmia. Sovellus on tehokas ja monipuolinen. Se on helposti siirrettävissä ja ylläpidettävissä nykyisissä laitteistoissa, koska edullisen suoritustemuodon mukaisesti käytetty Java-pohjainen DEL-suoritin on helposti integroitavissa eri systeemeihin. Edullisesti sovelluksen käyttöliittymä toimii esimerkiksi Windows 2000/NT-ympäristössä, johon on asennettu Microsoftin Internet Explorer-selainen versio 5 tai jokin uudempi ratkaisu. Lisäksi koodingenerointisovellus on mukautuva ja helposti uudelleen käytettävissä ihan uudenlaisenkin lähdemateriaalin kanssa.  
 25

## Patenttivaatimukset

1. Menetelmä tietyssä alkuperäisessä datassa sijaitsevan tiedon käsittelemiseksi, **tunnettu** siitä, että menetelmä sisältää vaiheet, joissa
  - osoitetaan vähintään kaksi esimerkkitapausta (301, 302),
- 5   – osoitettujen esimerkkitapausten (301, 302) perusteella muodostetaan säännöstö (103) halutunlaisen tiedon irrottamiseksi, ja
  - muodostetun säännösten mukaiset data-alueet alkuperäisestä tietueesta irrotetaan (104).
2. Patenttivaatimuksen 1 mukainen menetelmä, **tunnettu** siitä, että irrotetut data-  
10   alueet muunnetaan formaatiltaan yhtenäisiksi (105).
3. Patenttivaatimuksen 1 mukainen menetelmä, **tunnettu** siitä, että osoitetuilla vähintään kahdella esimerkkitapauksella on kullakin rakenne ja sisältö, jolloin rakenne on esimerkkitapausten välillä samanlainen mutta sisältö on erilainen.
4. Patenttivaatimuksen 1 mukainen menetelmä, **tunnettu** siitä, että esimerkkitapa-  
15   paukset osoitetaan alkuperäisestä datasta (101).
5. Patenttivaatimuksen 1 mukainen menetelmä, **tunnettu** siitä, että esimerkkitapausten (301, 302) perusteella muodostettu säännöstö (103) tallennetaan myöhem-  
pää käyttöä varten.
6. Patenttivaatimuksen 1 mukainen menetelmä, **tunnettu** siitä, että osoitetut esi-  
20   merkkitapaukset tokenisoidaan (201) ennen niiden varsinaista käsittelyä korvaamalla tietyt esimerkkitapausten osiot niitä vastaavilla tietorakenteilla, jotka sisältävät tun-  
nisteiden, kuten tyyppimerkinnän tai nimen, sekä kyseisen osion datasisällön.
7. Patenttivaatimuksen 6 mukainen menetelmä, **tunnettu** siitä, että osoitetuissa vähintään kahdessa esimerkkitapauksessa (301, 302) on vähintään yksi samanlainen  
25   osio, jota esimerkkitapausten käsittelyssä vastaa tietty tokeni.
8. Patenttivaatimuksen 6 mukainen menetelmä, **tunnettu** siitä, että säännösten muodostamiseksi
  - pisin osoitetuista tokenisoiduista esimerkeistä merkitään säännölliseksi lausek-  
keeksi R (203, 301),
  - 30   – seuraavaksi pisintä osoitetuista tokenisoiduista esimerkeistä merkitään esimerkki-  
lausekkeeksi E (204, 302) ja

– verrataan säännöllistä lauseketta R (203) ja sen hetkistä esimerkkilauseetta E (204).

9. Patenttivaatimuksen 8 mukainen menetelmä, **tunnettu** siitä, että säännöllistä lauseketta R ja sen hetkistä esimerkkilauseketta E vertaillaan tietyn vertailualgoritmin avulla (205), joka palauttaa editointiskriptin D (206, 303, 306).

10. Patenttivaatimuksen 9 mukainen menetelmä, **tunnettu** siitä, että säännöllistä lauseketta R ja sen hetkistä esimerkkilauseketta E vertaillaan sellaisen vertailualgoritmin avulla (205), joka palauttaa lyhimmän mahdollisen editointiskriptin.

11. Patenttivaatimuksen 9 mukainen menetelmä, **tunnettu** siitä, että säännösten muodostamiseksi säännöllistä lauseketta R muunnetaan editointiskriptin D (207) sisältämien muutostietojen pohjalta.

12. Patenttivaatimuksen 8 mukainen menetelmä, **tunnettu** siitä, että muodostettu säännöllinen lauseke muodostaa säännösten (211).

13. Patenttivaatimuksen 1 mukainen menetelmä, **tunnettu** siitä, että muodostetun säännösten avulla alkuperäisestä datasta irrotetaan esimerkkitapausten mukaiset osiot (104).

14. Laitteisto tietyssä alkuperäisessä datassa sijaitsevan tiedon käsittelemiseksi, **tunnettu** siitä, että laitteisto sisältää välineet

- vähintään kahden esimerkkitapausten (301, 302) osoittamiseksi alkuperäisestä datasta (101),
- säännösten muodostamiseksi osoitettujen esimerkkitapausten (301, 302) perusteella halutunlaisen tiedon irrottamiseksi, ja
- muodostetun säännösten mukaisten data-alueiden irrottamiseksi (104) alkuperäisestä lähdemateriaalista.

15. Patenttivaatimuksen 14 mukainen laitteisto, **tunnettu** siitä, että siinä on välineet irrotettujen osioiden muuntamiseksi formaatiltaan yhtenäisiksi (105).

16. Patenttivaatimuksen 14 mukainen laitteisto, **tunnettu** siitä, että esimerkkitapausten osoittamista varten laitteistossa on osoittimet merkkijonoihin.

17. Patenttivaatimuksen 14 mukainen laitteisto, **tunnettu** siitä, että laitteisto sisältää välineet osoitettujen esimerkkien tokenisoimiseksi korvaamalla tietyt esimerkkitapausten osiot niitä vastaavilla tietorakenteilla, jotka sisältävät tyyppimerkinnän tai nimen sekä kyseisen osion datasisällön.



18. Patenttivaatimuksen 17 mukainen laitteisto, **tunnettu** siitä, että laitteisto sisältää välineet tokenisoidun datan käsittelemiseksi.

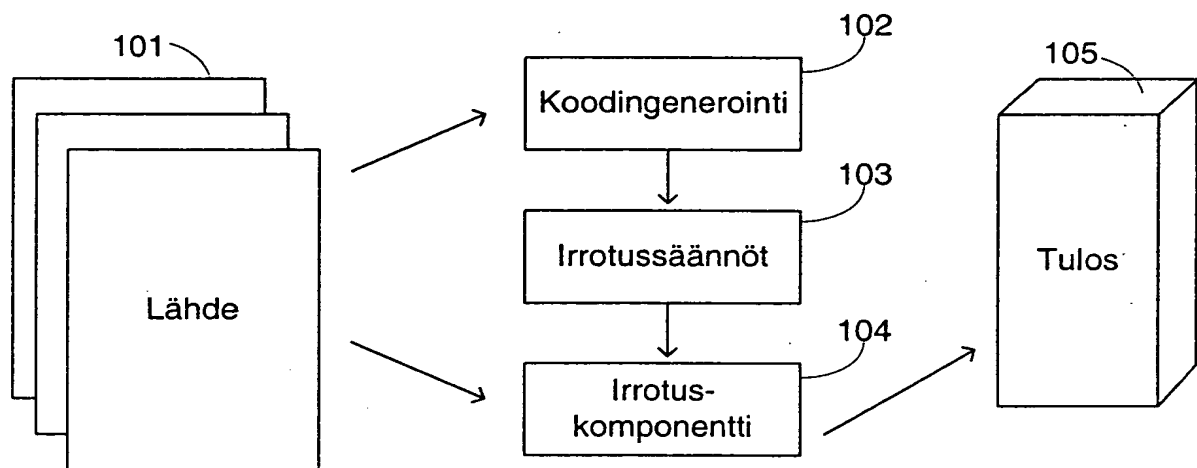
19. Patenttivaatimuksen 14 mukainen laitteisto, **tunnettu** siitä, että laitteisto sisältää välineet säännöstön muodostamiseksi muodostetun säännöllisen lausekkeen R mukaisesti.

20. Patenttivaatimuksen 19 mukainen laitteisto, **tunnettu** siitä, että välineet säännöstön muodostamiseksi sisältävät tätä tarkoitusta varten luodun ohjelmistokomponentin (102), joka on erillään sellaisesta ohjelmistokomponentista (104), joka on tarkoitettu data-alueiden irrottamiseen muodostettua säännöstöä käyttäen.

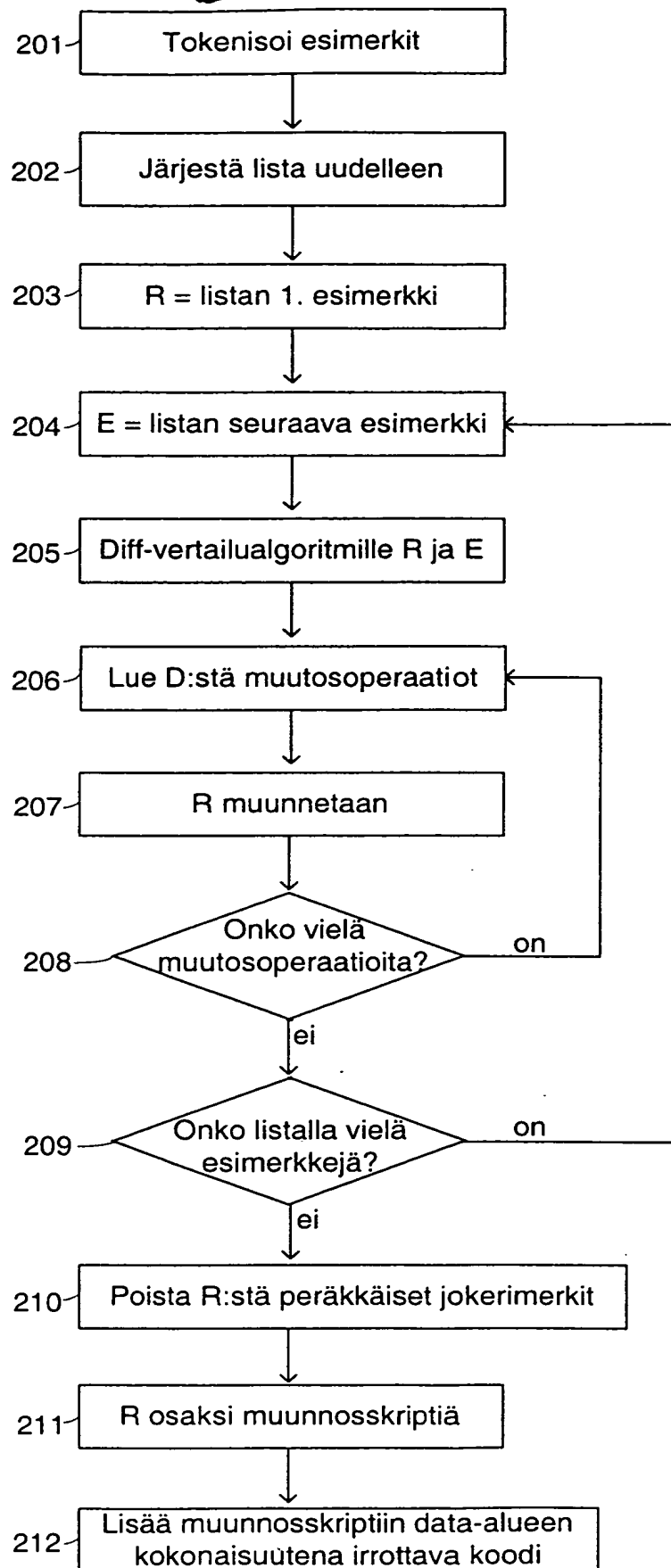
### (57) Tiivistelmä

Keksintö koskee sellaisten sääntöjen muodostamismenetelmää, jonka avulla tietoa voidaan ryhmitellä uudelleen. Lisäksi keksintö koskee järjestelyä tällaisen menetelmän toteuttamiseksi. Keksinnön tavoitteena on tuottaa menetelmä ja järjestely, joilla ohjelmointitaidotonkin käyttäjä voi muodostaa lähdedatasta valitsemilleen data-alueille irrotussäännöt. Kyseessä olevassa menetelmässä alkuperäisessä datassa olevaa tietoa käsitellään siten, että alkuperäisestä lähdedatasta (101) käyttäjä valitsee vähintään kaksi esimerkkitapausta (301, 302). Osoitettujen esimerkkitapausten perusteella muodostetaan säännöstö (103), jonka mukaiset data-alueet alkuperäisestä lähdedatasta irrotetaan (104). Irrotettuja data-alueita (105) voidaan jatkokäsitellä halutulla tavalla.

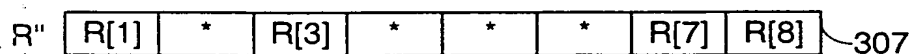
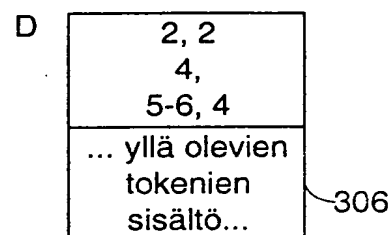
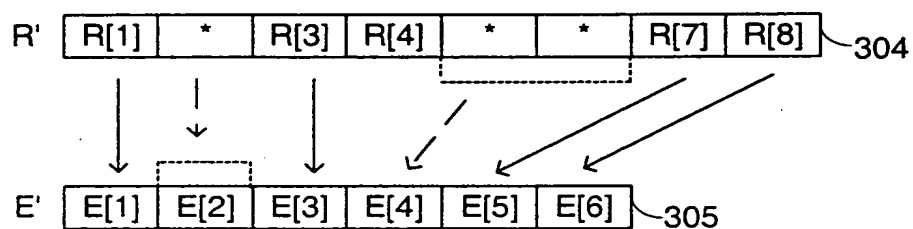
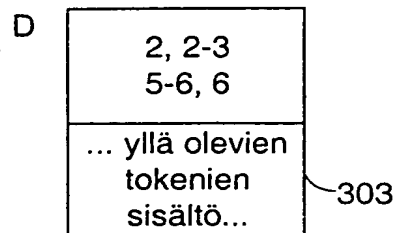
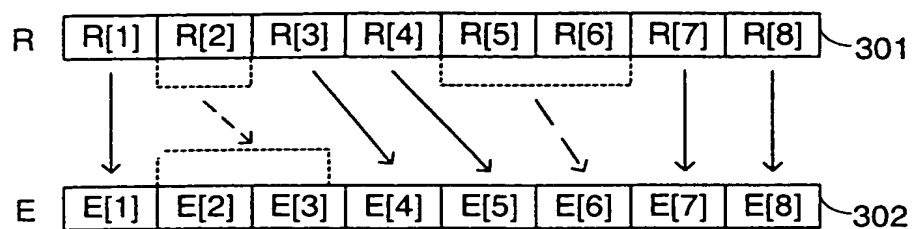
Kuva 1



Kuva 1.



Kuva 2.



Kuva 3.